1    **Title:** Signatures of selection at drug resistance loci in *Mycobacterium tuberculosis*

2    **Authors:** Tatum D. Mortimer[1], Alexandra M. Weber[1], Caitlin S. Pepperell[1*]

3    [1] Department of Medicine, Division of Infectious Diseases and Department of Medical
4    Microbiology and Immunology, University of Wisconsin-Madison

5    *Address correspondence to Caitlin S. Pepperell, cspepper@medicine.wisc.edu

6    Running title: Selection at drug resistance loci in *M. tb*

7

8    Word counts:

9    Abstract- 250

10    Importance- 134

11    Main Text- 5,560

12

## Abstract:

Tuberculosis (TB) is the leading cause of death by an infectious disease, and global TB control efforts are increasingly threatened by drug resistance in *Mycobacterium tuberculosis* (*M. tb*). Unlike most bacteria, where lateral gene transfer is an important mechanism of resistance acquisition, resistant *M. tb* arises solely by *de novo* chromosomal mutation. Using whole genome sequencing data from two natural populations of *M. tb*, we characterized the population genetics of known drug resistance loci using measures of diversity, population differentiation, and convergent evolution. We found resistant sub-populations to be less diverse than susceptible sub-populations, consistent with ongoing transmission of resistant *M. tb*. A subset of resistance genes ("sloppy targets") were characterized by high diversity and multiple rare variants; we posit that a large genetic target for resistance and relaxation of purifying selection contribute to high diversity at these loci. For "tight targets" of selection, the path to resistance appeared narrower, evidenced by single favored mutations that arose numerous times on the phylogeny and segregated at markedly different frequencies in resistant and susceptible sub-populations. These results suggest that diverse genetic architectures underlie drug resistance in *M. tb*, and combined approaches are needed to identify causal mutations. Extrapolating from patterns observed in well-characterized genes, we identified novel candidate variants involved in resistance. The approach outlined here can be extended to identify resistance variants for new drugs, to investigate the genetic architecture of resistance, and, when phenotypic data are available, to find candidate genetic loci underlying other positively selected traits in clonal bacteria.

## Importance:

*Mycobacterium tuberculosis* (*M. tb*), the causative agent of tuberculosis (TB), is a significant burden on global health. Antibiotic treatment imposes strong selective pressure on *M. tb* populations. Identifying the mutations that cause drug resistance in *M. tb* is important for guiding TB treatment and halting the spread of drug resistance. Whole genome sequencing (WGS) of *M. tb* isolates can be used to identify novel mutations mediating drug resistance and to predict resistance patterns faster than

43    traditional methods of drug susceptibility testing. We have used WGS from natural

44    populations of drug resistant *M. tb* to characterize effects of selection for advantageous

45    mutations on patterns of diversity at genes involved in drug resistance. The methods

46    developed here can be used to identify novel advantageous mutations, including new

47    resistance loci, in *M. tb* and other clonal pathogens.

48

**Introduction:**

*Mycobacterium tuberculosis* (*M. tb*), the causative agent of tuberculosis (TB), is estimated to have caused 1.4 million deaths in 2015, making it the leading cause of death due to an infectious disease. The proportion of TB due to MDR (multi-drug resistant) *M. tb* resistant to first line anti-tuberculosis drugs isoniazid (INH) and rifampin (RIF)) is increasing (1), which poses a major threat to global public health. Unlike most bacteria, *M. tb* evolves clonally, so resistance cannot be transferred among strains or acquired from other species of bacteria: drug resistance in *M. tb* results from *de novo* mutation within patients and transmission of drug resistant clones (2–4). The relative contributions of *de novo* emergence and transmitted drug resistance varies across sampling locations (5–9). Another potential variable in the emergence of drug resistant TB is *M.tb*'s lineage structure: seven distinct lineages have been identified among globally extant populations of *M. tb*. Among these, lineage 2 (L2) has been associated with relatively high rates of drug resistance, and it has been postulated that the acquisition of resistance is a result of higher rates of mutation in this lineage (10). Studies of *M.tb* evolution within hosts with TB have shown that emergence of drug resistance is associated with clonal replacements that lead to reductions in genetic diversity of the bacterial population (11, 12).

Many of the methods developed to identify advantageous mutations, such as those conferring antibiotic resistance, depend on recombination to differentiate target loci from neutral variants (13). However, in clonal organisms like *M. tb*, neutral and deleterious mutations that are linked to advantageous variants will evolve in tandem with them. This linkage among sites can also cause competition between genetic backgrounds with beneficial mutations, decreasing the rate of fixation for beneficial alleles, while deleterious alleles are purged less efficiently (14–16).

While the majority of the *M. tb* genome is subject to purifying selection (i.e. selection against deleterious mutations) (3), antibiotic pressure exerts strong selection for advantageous variants that confer resistance. *M. tb* drug resistance has been the focus of extensive investigation, and a variety of resistance mutations have been characterized for commonly used anti-tuberculosis drugs (17). Drug resistance

79    mutations can be associated with fitness costs (18–20), and compensatory mutations

80    that ameliorate these fitness costs have been identified in the context of rifampicin

81    resistance (21, 22). Resistance mutations found to have lower fitness costs *in vitro* - as

82    measured by competition assays - are found at higher frequencies among *M. tb* clinical

83    isolates and appear to be transmitted at higher rates relative to mutations with high *in*

84    *vitro* fitness costs (18, 23). Candidate loci involved in resistance and compensation for

85    its fitness effects have been identified previously by screening for homoplastic variants

86    (i.e. mutations that emerge more than once on the phylogeny) that are significantly

87    associated with drug resistant phenotypes (24) and genes with higher *dN/dS* (ratio of

88    non-synonymous versus synonymous mutations) in resistant compared to sensitive

89    isolates (25). Application of these methods to whole genome sequence data from *M. tb*

90    clinical isolates has recovered known drug resistance loci, as well as loci associated

91    with cell surface lipids and biosynthesis, DNA replication, and metabolism.

92    The goal of the present study was to use patterns of genetic diversity at known drug

93    resistance loci to identify the population genomic signatures of positive selection in

94    natural populations of *M. tb*. Using whole genome sequence data from two populations

95    for which phenotypic resistance data were available, we have identified several distinct

96    signatures associated with these loci under selection. Based on these results, we

97    propose methods of identifying loci under positive selection, including novel drug

98    resistance loci, in clonal bacteria such as *M. tb*.

99    **Results:**

100   We inferred the phylogeny of 1161 *M. tb* isolates from Russia and South Africa (see

101   Methods, Supplementary Table 1) using the approximate maximum likelihood method

102   implemented in FastTree2 (Figure 1). The majority of the isolates belong to L2 ($n = 667$)

103   and L4 ($n = 481$). *M. tb* nucleotide diversity was similar to previous estimates from a

104   globally distributed sample (26).  We identified lineage-specific patterns in overall

105   diversity, with L4 having higher diversity than L2 ($\pi_{L2}$: 3.6 x $10^{-5}$, $\pi_{L4}$: 1.5 x $10^{-4}$).

106   Previously published analyses of whole genome sequence data from L2 indicate that

107   the majority of L2 isolates worldwide belong to a sub-lineage that has undergone

108   relatively recent expansion (27, 28). In this sample from Russia and South Africa, the

109    majority of L2 isolates belong to this sub-lineage, while the L4 isolates are associated

110    with deeper branching sub-lineages. This likely contributes to the observed differences

111    in diversity.

112    Overall diversity of L2 was lower than L4 in our sample (Figure 2, $p < 2.2 \times 10^{-16}$).

113    Seven hundred and sixty two of the isolates in our sample (66%) are resistant to one or

114    more anti-tuberculosis drugs (Table 1).

115    Drug resistant TB can be acquired as a result of *de novo* mutations within a patient or

116    by infection with a resistant strain. When resistance is primarily mediated by *de novo*

117    mutations, diversity should be similar in susceptible and resistant populations as

118    resistance will arise on multiple genetic backgrounds. By contrast, if resistance

119    develops primarily *via* transmission of resistant strains, the resistant sub-population

120    should be less diverse than the susceptible sub-population. We compared the

121    nucleotide diversity of susceptible and resistant sub-populations and found  genome

122    wide estimates of nucleotide diversity to be higher in isolates susceptible to a range of

123    drugs for which phenotyping data were available (paired t-test, $p = 0.029$). In

124    comparisons of gene-wise diversity in susceptible and resistant populations, we found

125    that resistant isolates had a greater number of genes with no diversity, but levels of

126    diversity within genes in which it was measurable were similar between resistant and

127    susceptible populations (Figure 3).

128    Of the 3,162 genes included in our analyses, 109 (3%) were invariant across all isolates

129    in our sample. This is likely due to strong purifying selection on these genes. An

130    additional 136 genes harbored variation in the drug susceptible populations but were

131    invariant across all of the drug resistant populations. We did not observe the converse,

132    i.e. genes that were invariant in susceptible isolates specifically, which supports the

133    conclusion from genome wide diversity estimates that resistant isolates represent a

134    subset of the diversity found in susceptible populations and suggests that there may be

135    purifying selection that is specific to the setting of drug resistance. In order to evaluate

136    whether the observed pattern was likely to arise by chance, we performed weighted

137    random sampling of genes. The weighting was based on diversity in susceptible

138    populations, assuming that genes with low diversity in susceptible populations are more

139 likely to be invariant in resistant populations. After randomly sampling genes in each

140 drug resistant population 1000 times, we found that no samples contained shared

141 genes amongst all resistant populations (first and second line drugs). This suggests that

142 specific genes tend to lose diversity in the setting of drug resistance, which could result

143 from purifying selection specific to this setting. A potentially important caveat is that in

144 our data set, drug resistant populations are not independent and the same isolates are

145 often resistant to multiple drugs. Since resistance to second line drugs frequently arises

146 on genetic backgrounds already resistant to one or more first line drugs, we repeated

147 the sampling with only first line drugs and found that the maximum number of sampled

148 genes shared across all populations was 2 (median 0). Overall,  these results suggest

149 that certain genes are more likely to lose diversity as drug resistance evolves, but we

150 cannot completely rule out the possibility that the pattern arose as a result of

151 overlapping membership in resistant populations.

152 We compared diversity of drug resistance associated genes (Table 2) with the rest of

153 the genome using two measures of diversity: average pairwise differences ($\pi$) and

154 number of segregating sites ($\theta$). We found the resistance genes *gid*, *rpsL,* and *pncA* to

155 be in the top $5^{th}$ percentile of gene-wise $\pi$ and/or $\theta$ values. *rrs* and *ethA* are in the top

156 $5^{th}$ percentile of $\theta$*,* but not $\pi$.  Surprisingly, despite being a target of multiple drug

157 resistance mutations (Table 2), we did not identify extreme levels of diversity in *katG*

158 ($80^{th}$ and $82^{nd}$ percentile of $\pi$ and $\theta$, respectively).

159 We also examined gene-wise diversity values within each lineage to look for lineage

160 specific high diversity genes. In both L2 and L4, *gid, rpsL, pncA, ethA,* and *thyA* were in

161 the top $5^{th}$ percentile of diversity ($\pi$ and/or $\theta$).  In L2, *rpoB, embB, Rv1772*, and *folC*

162 were additionally in the top $5^{th}$ percentile of gene-wise $\pi$ and/or $\theta$ values. In L4, *Rv0340*

163 was in the top $5^{th}$ percentile of gene-wise $\pi$ and/or $\theta$. While *rpoB* and *embB* were not in

164 the top $5^{th}$ percentile of gene-wise $\theta$ in L4, they still had high diversity ($91^{st}$ and $82^{nd}$

165 percentile, respectively).  The lineage specific differences in diversity of *Rv1772*, *folC*,

166 and *Rv0340* suggest that there are interactions between these loci and loci that

167 differentiate L2 and L4.

168    We used gene-wise estimates of Tajima's D to investigate gene specific skews in the

169    site frequency spectrum that could result from selection, where negative values indicate

170    an excess of rare variants and positive values indicate an excess of intermediate

171    frequency variants. We previously identified a relationship between gene length and

172    gene-wise estimates of Tajima's D for *M. tb* (26), and this finding was corroborated here

173    ($R^2$ = 0.3 after $\log_2$ transformation). In order to identify genes with extreme values of

174    Tajima's D - out of proportion with their length - we performed linear regression on $\log_2$

175    transformed gene lengths and Tajima's D values and identified genes with the largest

176    residuals (Figure 4). *pncA, ethA*, and *embC* all had Tajima's D values lower than

177    expected based on their length ($5^{th}$ percentile of residual values). This indicates that

178    these genes contain an excess of rare variants compared to other genes in the genome.

179    Excess rare variants can result from a population expansion, a selective sweep, or

180    purifying selection.

181    We calculated the ratio of π and θ of resistance associated genes in isolates

182    susceptible and resistant to first line drugs and identified genes with markedly different

183    diversities in resistant and susceptible sub-populations (Figure 5A). Among resistance

184    genes in the top $5^{th}$ percentile of gene-wise π and θ overall, diversity of *pncA* and *ethA*

185    is relatively high among resistant isolates, whereas diversity of *gid* is similar in resistant

186    and susceptible populations. We also examined differences in this ratio between

187    isolates in L2 and L4 (Figure 5B). *Rv1772* and *embR* were more diverse in resistant

188    isolates in L2, and *kasA* and *tlyA* were more diverse in resistant isolates in L4.

189    We used $F_{ST}$ outlier analysis to identify single nucleotide polymorphisms (SNPs) and

190    indels that exhibited extreme differences in frequency between susceptible and resistant

191    populations. Our *a priori* expectation was that variants mediating resistance would be at

192    markedly higher frequency in the drug resistant sub-population and that drug targets

193    would be enriched among genes harboring variants with high $F_{ST}$. After removing SNPs

194    in regions corresponding to indels and variants at sites missing data for greater than 5%

195    of isolates, the highest $F_{ST}$ SNPs in comparisons of resistant and susceptible sub-

196    populations to first line drugs are in *katG* (2155168, $F_{ST}$ = 0.89, INH)*, rpoB* (761155, $F_{ST}$

197    = 0.72, RIF), and *rpsL* (781687, $F_{ST}$ = 0.37, streptomycin (STR)). These SNPs were also

198   $F_{ST}$ outliers in the lineage specific analyses. We used a randomization procedure to

199   assess the significance of observed $F_{ST}$ values and found the maximum $F_{ST}$ values after

200   randomly assigning resistant and susceptible designations to be 0.023 for INH, 0.019

201   for RIF, and 0.018 for STR. In addition to SNPs within known drug resistance

202   associated genes, we identified $F_{ST}$ outliers in genes that may be novel targets for drug

203   resistance (Table 3).

204   Homoplastic SNPs – i.e. SNPs that evolve more than once on a phylogeny – are

205   candidate loci under positive selection and have previously been used to identify

206   resistance associated mutations in *M. tb* (24). Of the 235 genes with homoplastic SNPs

207   that we identified in our sample, 13 are known to be associated with drug resistance

208   (Figure 6), and resistance genes were significantly enriched among genes with

209   homoplastic SNPs (Fisher's Exact Test, $p = 1.2 \times 10^{-4}$). *pncA* had the largest number of

210   homoplastic SNPs of any gene in the genome ($n = 27$ distinct SNPs that appear $> 1$ on

211   the phylogeny). The SNPs identified in $F_{ST}$ analysis were also identified as homoplastic

212   (Figure 6). Our results suggest that complementary approaches based on homoplasy

213   and $F_{ST}$ outlier analysis can be used to identify SNPs associated with a trait of interest

214   (in this case drug resistance). In addition to genic SNPs, we observed homoplastic

215   SNPs that are also $F_{ST}$ outliers in intergenic regions upstream of drug resistance

216   associated genes (Table 3). These are candidate resistance and compensatory

217   mutations with a regulatory mechanism of action.

218   In our analyses of indels, we controlled for the possibility that indels affecting the same

219   gene may not be called in exactly the same position by considering indels within the

220   same gene as identical. We identified four drug resistance associated genes with

221   homoplastic indels: *gid, ethA*, *rpoB,* and *pncA*. $F_{ST}$ values for the deletion in *gid* were in

222   the top 5$^{th}$ percentile for capreomycin (CAP), ethambutol (EMB), ethionamide (Et),

223   kanamycin (K), ofloxacin (OFL), and pyrazinamide (PZA) resistant populations, but,

224   interestingly, the deletion was not associated with STR resistance ($F_{ST} = 0.04$). Unlike

225   homoplastic SNPs, homoplastic indels were not significantly enriched for drug

226   resistance associated loci ($p = 1$).

227    We recovered 20 out of 40 known drug targets by identifying genes with extreme values

228    of diversity, homoplastic SNPs, or SNPs that are $F_{ST}$ outliers in comparisons of resistant

229    and susceptible subpopulations.  All genes with both extremely high diversity (top 5[th]

230    percentile) and homoplastic mutations were drug resistance associated (i.e. *gid, ethA,*

231    *pncA,* and *rpsL*). We identified 67 genes with high diversity and Tajima's D values more

232    negative than expected based on gene length; only two of these were associated with

233    drug resistance (i.e. *ethA* and *pncA*). Twenty out of 51 homoplastic SNPs that are also

234    $F_{ST}$ outliers fall within or upstream of known drug resistance associated genes. The

235    remaining SNPs may be false positives or novel drug resistance mutations.

## Discussion:

237    Highly virulent bacterial pathogens such as *M. tb, Yersinia pestis* (29), *Francisella*

238    *tularensis* (30), and *Mycobacterium ulcerans* (31) appear to evolve clonally, i.e. with

239    little to no evidence of lateral gene transfer.  It is important to identify advantageous

240    mutations in these and other organisms, as they are likely to be associated with

241    phenotypes such as drug resistance, heightened transmissibility, or host adaptation.

242    However, few methods are available for identifying loci under positive selection in the

243    setting of clonal evolution. We adopted an empirical approach to this problem and used

244    natural population data to characterize patterns of diversity at loci known to be under

245    positive selection in *M. tb.*

246    In this analysis of clinical isolates from settings with endemic drug resistance, we found

247    genome-wide diversity to be higher in susceptible *M. tb* sub-populations than in those

248    resistant to first- and second- line drugs (with the exception of protionamide (PRO) and

249    moxifloxacin (MOX) resistant populations). The observation of higher diversity in drug

250    susceptible populations is consistent with a significant role for transmitted resistance in

251    the propagation of drug resistant *M. tb.* A recent study of extensively drug resistant

252    (XDR) *M. tb* infection in South Africa concluded that XDR cases result primarily from

253    transmission of resistance, rather than *de novo* evolution of resistance mutations during

254    infection (9).  The primary studies for the sequence data analyzed here also identified

255    clusters of drug resistant isolates (5, 6), suggesting that resistant isolates were being

256    transmitted. Our results, along with these previously published observations, suggest

257    that the fitness of drug resistant isolates can be high enough to allow them to circulate

258    in endemic regions. As discussed below, the fitness effects of *M. tb* drug resistance

259    mutations appear to vary substantially; the finding of transmitted resistance in this and

260    other studies suggests that the fitness of isolates harboring low-cost mutations is

261    comparable to that of susceptible *M. tb*. The populations in our study have a high

262    burden of drug resistant TB, and the role of transmitted drug resistance may differ in

263    other settings.

264    An alternative – but not mutually exclusive – explanation for the observation of higher

265    diversity in susceptible populations is that drug resistant *M. tb* is under distinct

266    evolutionary constraints that reduce average genome-wide levels of diversity. In support

267    of this hypothesis, we identified a specific subset of genes that were invariant across

268    drug resistant populations. Interestingly, while average diversity was lower for resistant

269    sub-populations, the gene-wise diversity distributions had heavier tails, indicating there

270    were more genes with extreme levels of diversity.

271

272    We found the genetic architecture of resistance to vary among targets, and resistance-

273    associated genes tended to fall within categories that we term "sloppy", "tight", and

274    "hybrid" targets of selection (the latter has a combination of tight and sloppy features

275    and applies to *rpsL, embB,* and *rpoB*). "Sloppy" resistance genes are characterized by

276    high levels of diversity. Genes associated with PZA, EMB, Et, and STR resistance (i.e.

277    *pncA*, *gid, rpsL, rrs, ethA*) have high levels of diversity; some also had an excess of rare

278    variants (*pncA*, *ethA*, *embC*). The finding that these genes accumulate multiple,

279    individually rare mutations implies that there is a large target for resistance and/or

280    compensatory mutations within the gene: that is, resistance can result from multiple

281    different variants acting individually or in concert. In addition to its numerous rare

282    mutations, *pncA* also contains the highest number of homoplastic SNPs (27 SNPs

283    emerged more than once on the phylogeny) of any gene in the data set. Among the 62

284    non-synonymous *pncA* mutations in our dataset, 55 have been previously reported in

285    association with drug resistance (TB Drug Resistance Mutation Database (32)). The

286    newly described SNPs may mediate drug resistance or compensation for the fitness

287   effects of other variants. Relaxed purifying selection is likely to play a role in concert

288   with selection for diverse advantageous resistance mutations in the accumulation of

289   diversity in *pncA* and other sloppy targets. The fact that numerous mutations are

290   segregating in a natural population suggests that alterations to these genes are

291   generally associated with negligible fitness costs.  An *M. tb* strain harboring a deletion in

292   *pncA* conferring resistance to PZA was estimated to be endemic in Quebec by 1800,

293   long before the use of PZA for the treatment of TB (33–35). This supports the idea that

294   purifying selection on *pncA* is relatively weak, which would contribute to its exceedingly

295   high diversity and broaden the adaptive paths to resistance.

296   In contrast to *pncA*, *gid,* which is associated with low level STR resistance (36), does

297   not appear to have the signatures of a "sloppy" target for resistance despite its high

298   diversity. We identified just three homoplastic SNPs within *gid*, and previous studies

299   have found that STR resistant isolates do not encode the same *gid* mutations (37). This

300   could indicate that a multitude of mutations within *gid* confer resistance, but levels of

301   diversity in the gene were similar in resistant and susceptible isolates. Previous studies

302   of sequence polymorphism in *gid* have identified high diversity in this gene in both

303   resistant and susceptible isolates (37–39): *gid* appears to be subject to relaxed purifying

304   selection in the presence and absence of antibiotic pressure. Since *gid* mutations confer

305   low level resistance, it's also possible that mis-classification of resistance phenotypes

306   contributed to the lack of differentiation we and others have observed between

307   putatively STR resistant and susceptible sub-populations. In addition, mutations in *rpsL*,

308   which cause high level resistance, could mask the contribution of *gid* to STR resistance.

309   We found some drug targets to be highly diverse in resistant sub-populations of either

310   L2 or L4 (but not both), suggesting that resistance mutations in these genes interact

311   with the genetic background; the fitness effects of mutations in these genes could, for

312   example, vary on different genetic backgrounds. Lineage-specific $F_{ST}$ outliers are

313   another category of candidate locus with lineage dependent roles in drug resistance

314   (Table 3). Epistatic interactions between drug resistance mutations and *M. tb* lineage

315   have been reported previously: for example, specific mutations in the *inhA* promoter

316   have been associated with the L1 and *M. africanum* genetic backgrounds (40, 41).

317    In contrast to "sloppy" targets, we discovered individual homoplastic SNPs associated

318    with drug resistant sub-populations (i.e. with high $F_{ST}$) representing "tight" targets of

319    selection in genes conferring resistance to INH, RIF, and STR. Numerous resistance

320    mutations have been described in *katG*, *rpoB, rpsL, embB,* and *gyrA,* but we find drug

321    resistant sub-populations to be defined by a specific subset of mutations in these genes.

322    This suggests that certain mutations are strongly favored relative to others conferring

323    resistance to the same drugs when *M. tb* is in its natural environment. Antibiotic

324    resistance can impose fitness costs on *M. tb* during *in vitro* growth, with the range of

325    fitness costs varying among mutations, even within the same gene (18). Mutations can

326    also have different fitness effects depending on the genetic background, but the most fit

327    mutants were the same across *M. tb* lineages in a study of RIF resistance (18).

328    In our analyses, we found the dominant INH resistance mutation in *katG* to affect the

329    serine at position 315. This change reduces affinity to INH but preserves catalase

330    activity (42) and is associated with lower fitness costs than other *katG* mutants, both *in*

331    *vitro*  and in a mouse model (43, 44). This mutation was recently shown to precede

332    mutations conferring resistance to other drugs during accumulation of resistance in

333    evolution of multi-drug resistant *M. tb* (45). The dominant mutations we identified in

334    *rpoB* (codon 450) and *rpsL* (codon 43) have also been found to have lower fitness costs

335    *in vitro* compared to other mutations conferring resistance to RIF and STR in these

336    genes (18, 44, 46). These results suggest that many of the findings regarding the

337    relative fitness costs of *M. tb* resistance mutations *in vitro* and in animal models are

338    relevant to the pathogen's natural environment.

339    The fitness effects of mutations in *gyrA* (codon 94) and *embB* (codon 306) have not

340    been measured; based on our homoplasy and $F_{ST}$ results, we propose that they have

341    lower fitness costs than other mutations in these genes and that they represent "tight"

342    targets of selection. Mutations at *gyrA* codon 94 were previously found to be the most

343    prevalent in a survey of *gyrA* and *gyrB* mutations in fluoroquinolone resistant *M. tb*

344    clinical isolates (47). Interestingly, the mutation in *embB* codon 306 has been previously

345    associated with acquisition of multiple resistances (48), and we find that this position is

346    an $F_{ST}$ outlier for all first line drugs in L4. This mutation is not an $F_{ST}$ outlier in L2 (i.e top

347    5<sup>th</sup> percentile), with percentiles for $F_{ST}$ values ranging from 0.07-0.68 for first line drugs

348    in this lineage. These observations suggest that the genetic background affects

349    interactions among resistance mutations, and that *embB* 306 is important for acquisition

350    of multidrug resistance in L4 but not L2.

351    We searched for indels with the signature of a "tight" target, i.e. homoplastic mutations

352    segregating at markedly different frequencies in drug susceptible and resistant sub-

353    populations. Unlike the pattern observed with SNPs, genes associated with drug

354    resistance were not significantly enriched among those harboring homoplastic indels.

355    We identified one homoplastic indel that was also an $F_{ST}$ outlier - a deletion in *gid* that

356    causes a frameshift. Patterns of variation in *gid* are complex and suggest a role for

357    relaxation of purifying selection (i.e. in the accumulation of excess SNPs in both

358    resistant and susceptible isolates) and perhaps a tight target associated with multi-

359    resistance (i.e. this homoplastic/$F_{ST}$ outlier deletion that was associated with resistance

360    to CAP, EMB, Et, K, OFL, and PZA).

361    Our finding that, save for the frameshift mutation in *gid*, indels in resistance genes do

362    not have the signature of "tight" targets suggests that they are generally associated with

363    higher fitness costs than SNPs. Fifteen drug targets have been found in transposon

364    mutagenesis experiments to be essential for *M. tb* growth *in vitro*, including *rpoB* and

365    *rpsL*; deletions in these genes are likely to interrupt important functions (49). Deletions

366    in non-essential genes could also have fitness costs. Deletions in *katG,* which is non-

367    essential, can result in INH resistance but they are not observed as frequently in clinical

368    isolates as the KatG S315 SNP, particularly among transmitted INH-resistant strains

369    (23).

370    There are several limitations to our study. Resistance to multiple drugs was common in

371    our sample, and in some cases it was difficult to identify patterns of diversity and

372    population differentiation that were specific to individual drugs. Our results are also

373    limited by the accuracy with which drug resistance phenotypes were determined and a

374    lack of phenotypic data for some drugs (particularly second line drugs).  Our sample

375    was heavily skewed to lineages 2 and 4, and the results are not necessarily applicable

376    to other *M. tb* lineages. Finally, the data analyzed here were generated with short

377    sequencing read technologies, and we were thus limited to characterizing diversity in

378    regions of the *M. tb* genome that can be resolved with these methods: regions that were

379    masked from analysis (e.g. due to sequence repeats) may include unknown resistance

380    targets. We also used an L4 genome (H37Rv) as a reference, and gene content specific

381    to L2 may not have been identified.

382    We were not able to recover all drug resistance associated genes using the analyses

383    performed here. This is likely a result of limited phenotypic data for some drugs and

384    their associated targets (e.g. *thyA* and *folC*, which are associated with aminosalicylic

385    acid (PAS) resistance). Our list of drug targets was dominated by genes associated with

386    INH resistance, and signatures in genes that harbor rare resistance associated alleles

387    may be subtle compared to the KatG S315 mutation found at high frequency in drug

388    resistant populations.

389    We identified 31 SNPs that do not fall within the list of known drug resistance genes,

390    which both emerged more than once on the phylogeny (homoplasies) and were

391    segregating at markedly different frequencies in resistant and susceptible sub-

392    populations ($F_{ST}$ outliers). These SNPs may be novel resistance determinants; notably,

393    all non-synonymous SNPs within this group are in genes linked with drug resistance in

394    other studies (i.e. they are in genes encoding efflux pumps, genes differentially

395    regulated in resistant isolates or in response to the presence of drug, potential drug

396    targets, or genes in the same pathways as drug targets or resistance determinants)

397    (50–54). In addition to a direct, previously unrecognized role in resistance, these SNPs

398    could compensate for fitness costs of drug resistance. For example, we identified a

399    homoplastic $F_{ST}$ outlier in *rpoC*, and mutations in *rpoC* have been shown to compensate

400    for RIF resistance in experimental evolution studies (22).

401    Intriguingly, we found lipid metabolism genes to be enriched in the list of genes

402    harboring homoplastic SNPs ($p = 0.013$). We've previously shown that these genes

403    have extreme values of diversity in a global sample of *M. tb* isolates and within

404    individual hosts (26), suggesting that lipid metabolism genes may also be under positive

405    selection in *M. tb* populations. The results presented here could be extended by

406  phenotypic characterization of lipid profiles and identification of homoplastic variants

407  that are at markedly different frequencies in isolates with distinct lipid profiles.

408

409  Here we have used drug resistance loci in *M. tb* to identify the signatures of positive

410  selection in a clonal bacterium. We found these loci to be associated with distinct

411  patterns of diversity that likely reflect differing genetic architectures underlying the traits

412  under selection. The evolutionary path to resistance is broad for some drugs with

413  "sloppy targets", whereas for drugs with "tight targets" the means of acquiring resistance

414  appear more limited. This is likely due to fitness effects of resistance mutations in *M.*

415  *tb*'s natural environment, as numerous resistance mutations have been identified in tight

416  target genes. We also found evidence suggesting that there are important interactions

417  among loci during the evolution of resistance. Our results suggest that purifying

418  selection on a subset of genes intensifies in the setting of resistance, which could reflect

419  epistatic interactions and/or a response to the metabolic milieu imposed by

420  antimycobacterial agents. The results presented here can be used to create more

421  realistic models of resistance evolution in *M. tb* and to develop novel strategies of

422  preventing or mitigating the acquisition of resistance. For example, the narrow path to

423  resistance for drugs with tight targets reveals potentially exploitable vulnerabilities, as

424  does the finding of interdependencies among specific loci and the genetic background

425  in the evolution of resistance and multi-resistance. As new TB drugs become available

426  for clinical use, the approach outlined here can be extended to investigate their

427  architectures of resistance.

428  Efforts are underway to sequence and perform drug susceptibility testing on thousands

429  of *M. tb* isolates with the goal of creating an exhaustive catalogue of drug resistance

430  mutations and eventually using WGS to diagnose drug resistance in clinical settings

431  (CRyPTIC project, http://modmedmicro.nsms.ox.ac.uk/cryptic/, last accessed: May 24,

432  2017). We found that loci under positive selection can be identified using relatively

433  simple methods: "tight" targets are highly differentiated in their allele frequencies across

434  phenotypic groups (i.e. $F_{ST}$ outliers) and appear as homoplasies on the phylogeny;

435  "sloppy" targets are characterized by high diversity and/or low Tajima's D, as well as

436 homoplasies. Extrapolating from patterns observed among known resistance variants,

437 we have discovered new candidate regulatory and genic resistance variants. The

438 methods used in this study are widely available and should scale to analysis of the large

439 collections of genomic and phenotypic data that are currently being generated. This

440 approach can be extended to identify novel resistance loci in bacteria for which drug

441 susceptibility phenotypes are defined, as well as other positively selected loci in clonal

442 bacterial populations.

443 **Methods:**

444 Reference guided assembly

445 We downloaded sequencing read data from two large surveys of drug resistant *M. tb* in

446 Russia (5) and South Africa (6). We used FastQC (55) and TrimGalore (56) for quality

447 assessment and adaptor trimming of the reads. Trimmed reads were mapped to *M. tb*

448 H37Rv (NC_000962.3) using BWA-MEM v 0.7.12 (57). We used Samtools v 1.2 (58)

449 and Picard Tools (https://broadinstitute.github.io/picard/) for sorting, format conversion,

450 and addition of read group information. Variants were identified using Pilon v 1.16 (59).

451 A detailed description of the reference guided assembly pipeline is available at

452 https://github.com/pepperell-lab/RGAPepPipe. We removed isolates with mean

453 coverage less than 20X, isolates with percentage of the genome covered at 10X less

454 than 90%, isolates where a majority of reads did not map to H37Rv, and isolates where

455 greater than 10% of sites were unknown after mapping. The final data set contains 1161

456 *M. tb* isolates (Supplementary Table 1). The alignment was masked to remove repetitive

457 regions including PE/PPE genes.

458 Phylogenetic analysis

459 We estimated the approximately maximum likelihood phylogeny using the masked

460 alignment from reference guided assembly with FastTree-2.1.9 (60). We compiled

461 FastTree using the double precision option to accurately estimate branch lengths of

462 closely related isolates. We used FigTree (http://tree.bio.ed.ac.uk/software/figtree/) for

463 tree visualization.

464   SNP annotation

465   A VCF of single nucleotide variants was created from the masked alignment using SNP-

466   sites v 2.3.2 (61). SNPs were annotated using SnpEff v 4.1j  (62)  to identify

467   synonymous, non-synonymous, and intergenic SNPs based on the annotation of *M. tb*

468   H37Rv.

469   Indel identification

470   Insertions and deletions were identified during variant calling with Pilon. We used Emu

471   (63) to normalize indels across multiple isolates. We used a presence/absence matrix

472   for the normalized indels for further analyses of indel diversity.

473   Population genetics statistics

474   Whole genome and gene-wise diversity ($\pi$ and $\theta$) and neutrality (Tajima's D) statistics

475   were calculated using Egglib v 2.1.10 (64) for whole genome alignments and gene-wise

476   alignments. Isolates were further divided by lineage and drug resistance phenotype.

477   Sites with missing data due to indels or low quality base calls more than 5% of isolates

478   in the alignment were not included in calculation of statistics. Values of Tajima's D

479   showed a correlation with gene length in our sample. To find genes with extreme values

480   of Tajima's D, we performed linear regression in R (65) on log transformed Tajima's D

481   values and gene length and identified genes with large residual values. To  identify

482   alleles with marked differences in frequency in resistant and susceptible isolates, Weir

483   and Cockerham's $F_{ST}$ (66) was calculated using populations of resistant and susceptible

484   isolates for each drug using vcflib v1.0.0-rc0-262-g50a3 (https://github.com/vcflib/vcflib).

485   For non-biallelic SNPs, we calculated $F_{ST}$ for the two most common variants.

486   Homoplasy

487   We used TreeTime (67) to perform ancestral reconstruction and place SNPs and indels

488   on the phylogeny. We identified homoplastic SNPs and indels as those arising multiple

489   times on the phylogeny.

490   Data availability

491 Unless otherwise noted, all data and scripts associated with this study are available at

492 https://github.com/pepperell-lab/mtbDrugResistance.

493

507

## References:

1.  World Health Organization. 2016. Global tuberculosis report 2016.

2.  Supply P, Warren RM, Bañuls A-L, Lesjean S, Van Der Spuy GD, Lewis L-A, Tibayrenc M, Van Helden PD, Locht C. 2003. Linkage disequilibrium between minisatellite loci supports clonal evolution of Mycobacterium tuberculosis in a high tuberculosis incidence area. Mol Microbiol 47:529–538.

3.  Pepperell CS, Casto AM, Kitchen A, Granka JM, Cornejo OE, Holmes EC, Birren B, Galagan J, Feldman MW. 2013. The Role of Selection in Shaping Diversity of Natural M. tuberculosis Populations. PLoS Pathog 9:e1003543.

4.  Eldholm V, Balloux F. 2016. Antimicrobial Resistance in Mycobacterium tuberculosis: The Odd One Out. Trends Microbiol 24:637–648.

5.  Casali N, Nikolayevskyy V, Balabanova Y, Harris SR, Ignatyeva O, Kontsevaya I, Corander J, Bryant J, Parkhill J, Nejentsev S, Horstmann RD, Brown T, Drobniewski F. 2014. Evolution and transmission of drug-resistant tuberculosis in a Russian population. Nat Genet 46:279–286.

6.  Cohen KA, Abeel T, Manson McGuire A, Desjardins CA, Munsamy V, Shea TP, Walker BJ, Bantubani N, Almeida DV, Alvarado L, Chapman SB, Mvelase NR, Duffy EY, Fitzgerald MG, Govender P, Gujja S, Hamilton S, Howarth C, Larimer JD, Maharaj K, Pearson MD, Priest ME, Zeng Q, Padayatchi N, Grosset J, Young SK, Wortman J, Mlisana KP, O'Donnell MR, Birren BW, Bishai WR, Pym AS, Earl AM. 2015. Evolution of Extensively Drug-Resistant Tuberculosis over Four Decades: Whole Genome Sequencing and Dating Analysis of Mycobacterium tuberculosis Isolates from KwaZulu-Natal. PLoS Med 12:e1001880.

7.  Coscolla M, Barry PM, Oeltmann JE, Koshinsky H, Shaw T, Cilnis M, Posey J, Rose J, Weber T, Fofanov VY, Gagneux S, Kato-Maeda M, Metcalfe JZ. 2015. Genomic Epidemiology of Multidrug-resistant Mycobacterium tuberculosis During Transcontinental Spread. J Infect Dis jiv025.

8.  Guerra-Assunção JA, Crampin AC, Houben R, Mzembe T, Mallard K, Coll F, Khan P, Banda L, Chiwaya A, Pereira RPA, McNerney R, Fine PEM, Parkhill J, Clark TG, Glynn JR. 2015. Large-scale whole genome sequencing of M. tuberculosis provides insights into transmission in a high prevalence area. eLife 4:e05166.

9.  Shah NS, Auld SC, Brust JCM, Mathema B, Ismail N, Moodley P, Mlisana K, Allana S, Campbell A, Mthiyane T, Morris N, Mpangase P, van der Meulen H, Omar SV, Brown TS, Narechania A, Shaskina E, Kapwata T, Kreiswirth B, Gandhi NR. 2017. Transmission of Extensively Drug-Resistant Tuberculosis in South Africa. N Engl J Med 376:243–253.

10. Ford CB, Shah RR, Maeda MK, Gagneux S, Murray MB, Cohen T, Johnston JC, Gardy J, Lipsitch M, Fortune SM. 2013. Mycobacterium tuberculosis mutation rate estimates from different lineages predict substantial differences in the emergence of drug-resistant tuberculosis. Nat Genet 45:784–790.

547   11.  Eldholm V, Norheim G, von der Lippe B, Kinander W, Dahle UR, Caugant DA, Mannsåker
548         T, Mengshoel AT, Dyrhol-Riise AM, Balloux F. 2014. Evolution of extensively drug-
549         resistant Mycobacterium tuberculosisfrom a susceptible ancestor in a single patient.
550         Genome Biol 15:490.

551   12.  Black PA, Vos M de, Louw GE, Merwe R van der, Dippenaar A, Streicher EM, Abdallah
552         AM, Sampson SL, Victor TC, Dolby T, Simpson JA, Helden P van, Warren RM, Pain A.
553         2015. Whole genome sequencing reveals genomic heterogeneity and antibiotic purification
554         in Mycobacterium tuberculosis isolates. BMC Genomics 16:857.

555   13.  Vitti JJ, Grossman SR, Sabeti PC. 2013. Detecting Natural Selection in Genomic Data.
556         Annu Rev Genet 47:97–120.

557   14.  Hill WG, Robertson A. 1966. The effect of linkage on limits to artificial selection. Genet
558         Res.

559   15.  Felsenstein J. 1974. The Evolutionary Advantage of Recombination. Genetics 78:737–756.

560   16.  Birky CW, Walsh JB. 1988. Effects of linkage on rates of molecular evolution. Proc Natl
561         Acad Sci 85:6414–6418.

562   17.  Silva PEAD, Palomino JC. 2011. Molecular basis and mechanisms of drug resistance in
563         Mycobacterium tuberculosis: classical and new drugs. J Antimicrob Chemother 66:1417–
564         1430.

565   18.  Gagneux S, Long CD, Small PM, Van T, Schoolnik GK, Bohannan BJM. 2006. The
566         Competitive Cost of Antibiotic Resistance in Mycobacterium tuberculosis. Science
567         312:1944–1946.

568   19.  Böttger EC, Springer B. 2008. Tuberculosis: drug resistance, fitness, and strategies for
569         global control. Eur J Pediatr 167:141–148.

570   20.  Strauss OJ, Warren RM, Jordaan A, Streicher EM, Hanekom M, Falmer AA, Albert H,
571         Trollip A, Hoosain E, Helden PD van, Victor TC. 2008. Spread of a Low-Fitness Drug-
572         Resistant Mycobacterium tuberculosis Strain in a Setting of High Human
573         Immunodeficiency Virus Prevalence. J Clin Microbiol 46:1514–1516.

574   21.  Brandis G, Wrande M, Liljas L, Hughes D. 2012. Fitness-compensatory mutations in
575         rifampicin-resistant RNA polymerase. Mol Microbiol 85:142–151.

576   22.  Comas I, Borrell S, Roetzer A, Rose G, Malla B, Kato-Maeda M, Galagan J, Niemann S,
577         Gagneux S. 2012. Whole-genome sequencing of rifampicin-resistant Mycobacterium
578         tuberculosis strains identifies compensatory mutations in RNA polymerase genes. Nat
579         Genet 44:106–110.

580   23.  Gagneux S, Burgos MV, DeRiemer K, Enciso A, Muñoz S, Hopewell PC, Small PM, Pym
581         AS. 2006. Impact of Bacterial Genetics on the Transmission of Isoniazid-Resistant
582         Mycobacterium tuberculosis. PLOS Pathog 2:e61.

583   24.  Farhat MR, Shapiro BJ, Kieser KJ, Sultana R, Jacobson KR, Victor TC, Warren RM,
584         Streicher EM, Calver A, Sloutsky A, Kaur D, Posey JE, Plikaytis B, Oggioni MR, Gardy JL,

585    Johnston JC, Rodrigues M, Tang PKC, Kato-Maeda M, Borowsky ML, Muddukrishna B,
586    Kreiswirth BN, Kurepina N, Galagan J, Gagneux S, Birren B, Rubin EJ, Lander ES, Sabeti
587    PC, Murray M. 2013. Genomic analysis identifies targets of convergent positive selection in
588    drug-resistant Mycobacterium tuberculosis. Nat Genet advance online publication.

589    25.    Zhang H, Li D, Zhao L, Fleming J, Lin N, Wang T, Liu Z, Li C, Galwey N, Deng J, Zhou Y,
590    Zhu Y, Gao Y, Wang T, Wang S, Huang Y, Wang M, Zhong Q, Zhou L, Chen T, Zhou J,
591    Yang R, Zhu G, Hang H, Zhang J, Li F, Wan K, Wang J, Zhang X-E, Bi L. 2013. Genome
592    sequencing of 161 Mycobacterium tuberculosis isolates from China identifies genes and
593    intergenic regions associated with drug resistance. Nat Genet advance online publication.

594    26.    O'Neill MB, Mortimer TD, Pepperell CS. 2015. Diversity of Mycobacterium tuberculosis
595    across Evolutionary Scales. PLoS Pathog 11:e1005257.

596    27.    Luo T, Comas I, Luo D, Lu B, Wu J, Wei L, Yang C, Liu Q, Gan M, Sun G, Shen X, Liu F,
597    Gagneux S, Mei J, Lan R, Wan K, Gao Q. 2015. Southern East Asian origin and
598    coexpansion of Mycobacterium tuberculosis Beijing family with Han Chinese. Proc Natl
599    Acad Sci 112:8136–8141.

600    28.    O'Neill MB, Kitchen A, Zarley A, Aylward W, Eldholm V, Pepperell CS. 2017. Lineage
601    specific histories of Mycobacterium tuberculosis dispersal in Africa and Eurasia. bioRxiv
602    210161.

603    29.    Achtman M, Zurth K, Morelli G, Torrea G, Guiyoule A, Carniel E. 1999. Yersinia pestis, the
604    cause of plague, is a recently emerged clone of Yersinia pseudotuberculosis. Proc Natl
605    Acad Sci 96:14043–14048.

606    30.    Larsson P, Elfsmark D, Svensson K, Wikström P, Forsman M, Brettin T, Keim P,
607    Johansson A. 2009. Molecular Evolutionary Consequences of Niche Restriction in
608    Francisella tularensis, a Facultative Intracellular Pathogen. PLoS Pathog 5:e1000472.

609    31.    Vandelannoote K, Meehan CJ, Eddyani M, Affolabi D, Phanzu DM, Eyangoh S, Jordaens
610    K, Portaels F, Mangas K, Seemann T, Marsollier L, Marion E, Chauty A, Landier J,
611    Fontanet A, Leirs H, Stinear TP, Jong D, C B. 2017. Multiple Introductions and Recent
612    Spread of the Emerging Human Pathogen Mycobacterium ulcerans across Africa. Genome
613    Biol Evol 9:414–426.

614    32.    Sandgren A, Strong M, Muthukrishnan P, Weiner BK, Church GM, Murray MB. 2009.
615    Tuberculosis Drug Resistance Mutation Database. PLOS Med 6:e1000002.

616    33.    Nguyen D, Brassard P, Westley J, Thibert L, Proulx M, Henry K, Schwartzman K, Menzies
617    D, Behr MA. 2003. Widespread Pyrazinamide-Resistant Mycobacterium tuberculosis
618    Family in a Low-Incidence Setting. J Clin Microbiol 41:2878–2883.

619    34.    Nguyen D, Brassard P, Menzies D, Thibert L, Warren R, Mostowy S, Behr M. 2004.
620    Genomic Characterization of an Endemic Mycobacterium tuberculosis Strain: Evolutionary
621    and Epidemiologic Implications. J Clin Microbiol 42:2573–2580.

622    35.    Brassard P, Henry KA, Schwartzman K, Jomphe M, Olson SH. 2008. Geography and
623    genealogy of the human host harbouring a distinctive drug-resistant strain of tuberculosis.
624    Infect Genet Evol 8:247–257.

625  36.  Okamoto S, Tamaru A, Nakajima C, Nishimura K, Tanaka Y, Tokuyama S, Suzuki Y, Ochi
626       K. 2007. Loss of a conserved 7-methylguanosine modification in 16S rRNA confers low-
627       level streptomycin resistance in bacteria. Mol Microbiol 63:1096–1106.

628  37.  Spies FS, Ribeiro AW, Ramos DF, Ribeiro MO, Martin A, Palomino JC, Rossetti MLR,
629       Silva PEA da, Zaha A. 2011. Streptomycin Resistance and Lineage-Specific
630       Polymorphisms in Mycobacterium tuberculosis gidB Gene. J Clin Microbiol 49:2625–2630.

631  38.  Feuerriegel S, Oberhauser B, George AG, Dafae F, Richter E, Rüsch-Gerdes S, Niemann
632       S. 2012. Sequence analysis for detection of first-line drug resistance in Mycobacterium
633       tuberculosis strains from a high-incidence setting. BMC Microbiol 12:90.

634  39.  Jagielski T, Ignatowska H, Bakuła Z, Dziewit Ł, Napiórkowska A, Augustynowicz-Kopeć E,
635       Zwolska Z, Bielecki J. 2014. Screening for Streptomycin Resistance-Conferring Mutations
636       in Mycobacterium tuberculosis Clinical Isolates from Poland. PLOS ONE 9:e100078.

637  40.  Homolka S, Meyer CG, Hillemann D, Owusu-Dabo E, Adjei O, Horstmann RD, Browne
638       ENL, Chinbuah A, Osei I, Gyapong J, Kubica T, Ruesch-Gerdes S, Niemann S. 2010.
639       Unequal distribution of resistance-conferring mutations among Mycobacterium tuberculosis
640       and Mycobacterium africanum strains from Ghana. Int J Med Microbiol 300:489–495.

641  41.  Fenner L, Egger M, Bodmer T, Altpeter E, Zwahlen M, Jaton K, Pfyffer GE, Borrell S,
642       Dubuis O, Bruderer T, Siegrist HH, Furrer H, Calmy A, Fehr J, Stalder JM, Ninet B, Böttger
643       EC, Gagneux S, Group  for the SHCS and the SME of TS. 2012. Effect of Mutation and
644       Genetic Background on Drug Resistance in Mycobacterium tuberculosis. Antimicrob
645       Agents Chemother 56:3047–3053.

646  42.  Wengenack NL, Uhl JR, Amand S, L A, Tomlinson AJ, Benson LM, Naylor S, Kline BC,
647       Cockerill FR, Rusnak F. 1997. Recombinant Mycobacterium tuberculosis KatG(S315T) Is a
648       Competent Catalase-Peroxidase with Reduced Activity toward Isoniazid. J Infect Dis
649       176:722–727.

650  43.  Pym AS, Saint-Joanis B, Cole ST. 2002. Effect of katG Mutations on the Virulence of
651       Mycobacterium tuberculosis and the Implication for Transmission in Humans. Infect Immun
652       70:4955–4960.

653  44.  Spies FS, von Groll A, Ribeiro AW, Ramos DF, Ribeiro MO, Dalla Costa ER, Martin A,
654       Palomino JC, Rossetti ML, Zaha A, da Silva PEA. 2013. Biological cost in Mycobacterium
655       tuberculosis with mutations in the rpsL, rrs, rpoB, and katG genes. Tuberculosis 93:150–
656       154.

657  45.  Manson AL, Cohen KA, Abeel T, Desjardins CA, Armstrong DT, Barry Iii CE, Brand J,
658       TBResist Global Genome Consortium, Chapman SB, Cho S-N, Gabrielian A, Gomez J,
659       Jodals AM, Joloba M, Jureen P, Lee JS, Malinga L, Maiga M, Nordenberg D, Noroc E,
660       Romancenco E, Salazar A, Ssengooba W, Velayati AA, Winglee K, Zalutskaya A, Via LE,
661       Cassell GH, Dorman SE, Ellner J, Farnia P, Galagan JE, Rosenthal A, Crudu V,
662       Homordean D, Hsueh P-R, Narayanan S, Pym AS, Skrahina A, Swaminathan S, Van der
663       Walt M, Alland D, Bishai WR, Cohen T, Hoffner S, Birren BW, Earl AM. 2017. Genomic
664       analysis of globally diverse Mycobacterium tuberculosis strains provides insights into the
665       emergence and spread of multidrug resistance. Nat Genet advance online publication.

666    46.    Billington OJ, McHugh TD, Gillespie SH. 1999. Physiological Cost of Rifampin Resistance
667            Induced In Vitro in Mycobacterium tuberculosis. Antimicrob Agents Chemother 43:1866–
668            1869.

669    47.    Eilertson B, Maruri F, Blackman A, Herrera M, Samuels DC, Sterling TR. 2014. High
670            Proportion of Heteroresistance in gyrA and gyrB in Fluoroquinolone-Resistant
671            Mycobacterium tuberculosis Clinical Isolates. Antimicrob Agents Chemother 58:3270–
672            3275.

673    48.    Hazbón MH, Valle MB del, Guerrero MI, Varma-Basil M, Filliol I, Cavatore M, Colangeli R,
674            Safi H, Billman-Jacobe H, Lavender C, Fyfe J, García-García L, Davidow A, Brimacombe
675            M, León CI, Porras T, Bose M, Chaves F, Eisenach KD, Sifuentes-Osornio J, León AP de,
676            Cave MD, Alland D. 2005. Role of embB Codon 306 Mutations in Mycobacterium
677            tuberculosis Revisited: a Novel Association with Broad Drug Resistance and IS6110
678            Clustering Rather than Ethambutol Resistance. Antimicrob Agents Chemother 49:3794–
679            3802.

680    49.    DeJesus MA, Gerrick ER, Xu W, Park SW, Long JE, Boutte CC, Rubin EJ, Schnappinger
681            D, Ehrt S, Fortune SM, Sassetti CM, Ioerger TR. 2017. Comprehensive Essentiality
682            Analysis of the Mycobacterium tuberculosis Genome via Saturating Transposon
683            Mutagenesis. mBio 8:e02133-16.

684    50.    Shekar S, Yeo ZX, Wong JCL, Chan MKL, Ong DCT, Tongyoo P, Wong S-Y, Lee ASG.
685            2014. Detecting Novel Genetic Variants Associated with Isoniazid-Resistant
686            Mycobacterium tuberculosis. PLOS ONE 9:e102383.

687    51.    Danilchanka O, Mailaender C, Niederweis M. 2008. Identification of a Novel Multidrug
688            Efflux Pump of Mycobacterium tuberculosis. Antimicrob Agents Chemother 52:2503–2511.

689    52.    Fu LM, Shinnick TM. 2007. Genome-wide exploration of the drug action of capreomycin on
690            Mycobacterium tuberculosis using Affymetrix oligonucleotide GeneChips. J Infect 54:277–
691            284.

692    53.    Phong TQ, Ha DTT, Volker U, Hammer E. 2015. Using a Label Free Quantitative
693            Proteomics Approach to Identify Changes in Protein Abundance in Multidrug-Resistant
694            Mycobacterium tuberculosis. Indian J Microbiol 55:219–230.

695    54.    Phelan J, Coll F, McNerney R, Ascher DB, Pires DEV, Furnham N, Coeck N, Hill-
696            Cawthorne GA, Nair MB, Mallard K, Ramsay A, Campino S, Hibberd ML, Pain A, Rigouts
697            L, Clark TG. 2016. Mycobacterium tuberculosis whole genome sequencing and protein
698            structure modelling provides insights into anti-tuberculosis drug resistance. BMC Med
699            14:31.

700    55.    Andrews S. 2012. FastQC.

701    56.    Kreuger F. 2013. TrimGalore!

702    57.    Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-
703            MEM. ArXiv13033997 Q-Bio.

704  58.  Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin
705      R. 2009. The Sequence Alignment/Map format and SAMtools. Bioinformatics 25:2078–
706      2079.

707  59.  Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
708      Wortman J, Young SK, Earl AM. 2014. Pilon: An Integrated Tool for Comprehensive
709      Microbial Variant Detection and Genome Assembly Improvement. PLOS ONE 9:e112963.

710  60.  Price MN, Dehal PS, Arkin AP. 2010. FastTree 2 – Approximately Maximum-Likelihood
711      Trees for Large Alignments. PLOS ONE 5:e9490.

712  61.  Page AJ, Taylor B, Delaney AJ, Soares J, Seemann T, Keane JA, Harris SR. 2016. SNP-
713      sites: rapid efficient extraction of SNPs from multi-FASTA alignments. biorxiv;038190v1.

714  62.  Cingolani P, Platts A, Wang LL, Coon M, Nguyen T, Wang L, Land SJ, Lu X, Ruden DM.
715      2012. A program for annotating and predicting the effects of single nucleotide
716      polymorphisms, SnpEff. Fly (Austin) 6:80–92.

717  63.  Salazar A, Earl A, Desjardins C, Abeel T. 2015. Normalizing alternate representations of
718      large sequence variants across multiple bacterial genomes. BMC Bioinformatics 16:A8.

719  64.  De Mita S, Siol M. 2012. EggLib: processing, analysis and simulation tools for population
720      genetics and genomics. BMC Genet 13:27.

721  65.  R Core Team. 2015. R: A Language and Environment for Statistical Computing. R
722      Foundation for Statistical Computing, Vienna, Austria.

723  66.  Weir BS, Cockerham CC. 1984. Estimating F-Statistics for the Analysis of Population
724      Structure. Evolution 38:1358–1370.

725  67.  Sagulenko P, Puller V, Neher R. 2017. TreeTime: maximum likelihood phylodynamic
726      analysis. bioRxiv 153494.

727
728

**Table 1. Frequency of resistance in data set.** AMI- amikacin, CAP- capreomycin, EMB- ethambutol, Et- ethionamide, INH- isoniazid, K- kanamycin, MOX- moxifloxacin, OFL- ofloxacin, PRO- protionamide, PZA- pyrazinamide, RIF- rifampin, STR- streptomycin.

| Drug | Resistant | Susceptible | Unknown |
|------|-----------|-------------|---------|
| INH | 0.59 | 0.33 | 0.08 |
| STR | 0.53 | 0.39 | 0.07 |
| RIF | 0.50 | 0.43 | 0.07 |
| EMB | 0.30 | 0.59 | 0.11 |
| PZA | 0.21 | 0.67 | 0.12 |
| OFL | 0.16 | 0.39 | 0.45 |
| PRO | 0.15 | 0.22 | 0.62 |
| CAP | 0.10 | 0.41 | 0.49 |
| MOX | 0.09 | 0.41 | 0.49 |
| Et | 0.06 | 0.07 | 0.88 |
| AMI | 0.05 | 0.34 | 0.61 |
| K | 0.05 | 0.12 | 0.83 |

**Table 2. Signatures of selection in known drug resistance genes.** The number of distinct entries in the TB Drug Resistance Mutation Database for each gene is reported in TB Dream column. π and θ are the percentiles for each diversity value, respectively. TD is the percentile of the residual after linear regression of Tajima's D with gene length. Genes with homoplastic SNPs are indicated with 'Y' in the Homoplasy column. If a homoplastic SNP was also an $F_{ST}$ outlier, it is indicated with a 'Y' in the $F_{ST}$ column. Genes are classified as tight, sloppy, or hybrid targets of selection based on diversity, homoplasy, and $F_{ST}$ results. (IG) indicates an intergenic SNP.

| Gene | Rv Number | Drug | TB Dream | π | θ | TD | Homoplasy | $F_{ST}$ | Type |
|------|-----------|------|----------|---|---|----|-----------|----------|------|
| katG | Rv1908c | INH | 226 | 0.80 | 0.82 | 0.34 | Y | Y | tight |
| pncA | Rv2043c | PZA | 195 | 0.97 | 1.00 | 0.00 | Y | N | sloppy |
| embB | Rv3795 | EMB | 117 | 0.77 | 0.89 | 0.08 | Y | Y | hybrid |
| ahpC | Rv2428 | INH | 31 | 0.20 | 0.21 | 0.61 | Y | N | - |
| tlyA | Rv1694 | CAP | 28 | 0.37 | 0.89 | 0.06 | N | N | - |
| embC | Rv3793 | EMB | 28 | 0.59 | 0.74 | 0.01 | N | N | - |
| embR | Rv1267c | EMB | 25 | 0.46 | 0.49 | 0.28 | N | N | - |
| rrs | Rvnr01 | STR, K, CAP | 24 | 0.89 | 1.00 | 0.08 | N | N | - |
| ethA | Rv3854c | Et | 23 | 0.72 | 1.00 | 0.00 | Y | Y (IG) | sloppy, tight (IG) |
| gid | Rv3919c | STR | 22 | 1.00 | 1.00 | 0.07 | Y | N | sloppy |
| gyrB | Rv0005 | MOX, OFL | 15 | 0.58 | 0.91 | 0.00 | Y | N | - |
| fabG1 | Rv1483 | INH, Et | 13 | 0.60 | 0.66 | 0.30 | Y | Y (IG) | tight |
| inhA | Rv1484 | INH, Et | 13 | 0.56 | 0.59 | 0.32 | Y | N | - |
| rpsL | Rv0682 | STR | 13 | 0.99 | 0.95 | 0.80 | Y | Y | hybrid |
| gyrA | Rv0006 | MOX, OFL | 12 | 0.81 | 0.94 | 0.10 | Y | Y | tight |
| embA | Rv3794 | EMB | 11 | 0.77 | 0.38 | 0.79 | N | N | - |
| kasA | Rv2245 | INH | 7 | 0.73 | 0.18 | 0.86 | N | N | - |
| ndh | Rv1854c | INH | 5 | 0.57 | 0.52 | 0.28 | N | N | - |
| iniA | Rv0342 | EMB, INH | 4 | 0.64 | 0.33 | 0.56 | N | N | - |
| Rv0340 | Rv0340 | INH | 3 | 0.89 | 0.88 | 0.57 | N | N | - |
| iniB | Rv0341 | EMB, INH | 3 | 0.07 | 0.07 | 0.79 | N | N | - |
| fbpC | Rv0129c | INH | 3 | 0.78 | 0.19 | 0.89 | N | N | - |
| rmlD | Rv3266c | EMB | 2 | 0.75 | 0.36 | 0.75 | N | N | - |
| iniC | Rv0343 | EMB, INH | 2 | 0.49 | 0.67 | 0.08 | N | N | - |
| thyA | Rv2764c | PAS | 2 | 0.84 | 0.94 | 0.28 | N | N | - |
| nat | Rv3566c | INH | 2 | 0.76 | 0.55 | 0.63 | N | N | - |
| accD6 | Rv2247 | INH | 1 | 0.90 | 0.63 | 0.90 | N | N | - |
| furA | Rv1909c | INH | 1 | 0.80 | 0.63 | 0.62 | N | N | - |
| Rv1772 | Rv1772 | INH | 1 | 0.50 | 0.35 | 0.54 | N | N | - |
| fabD | Rv2243 | INH | 1 | 0.26 | 0.28 | 0.54 | N | N | - |
| fadE24 | Rv3139 | INH | 1 | 0.36 | 0.58 | 0.12 | N | N | - |
| rpoB | Rv0667 | RIF | 1 | 0.82 | 0.92 | 0.18 | Y | Y | hybrid |
| efpA | Rv2846c | INH | 1 | 0.10 | 0.11 | 0.65 | N | N | - |

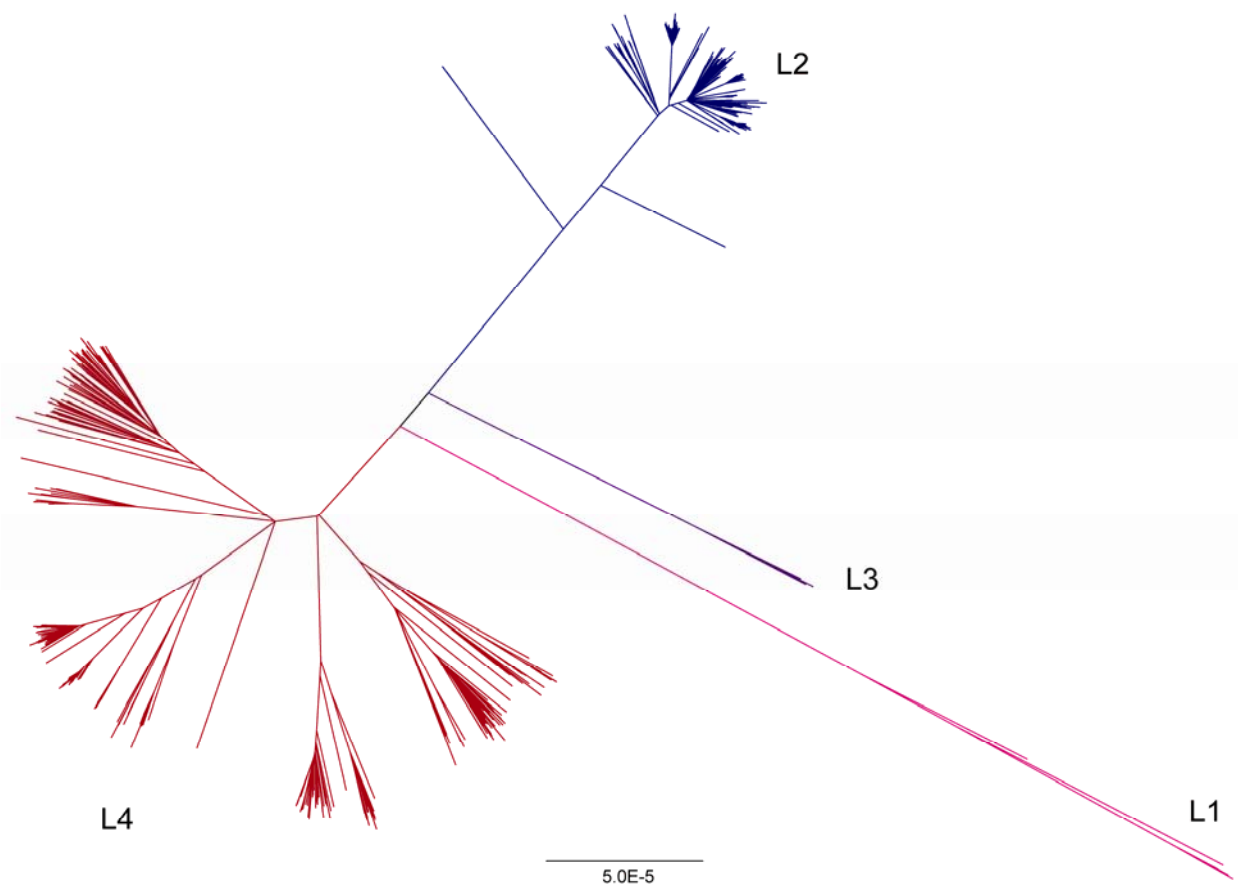| | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|
| *ethR* | *Rv3855* | Et | - | 0.58 | 0.77 | 0.22 | N | N | - |
| *Rv0678* | *Rv0678* | BDQ | - | 0.37 | 0.72 | 0.25 | N | N | - |
| *eis* | *Rv2416c* | K | - | 0.51 | 0.28 | 0.54 | N | N | - |
| *mshA* | *Rv0486* | Et | - | 0.86 | 0.48 | 0.87 | N | N | - |
| *rpsA* | *Rv1630* | PZA | - | 0.88 | 0.62 | 0.84 | N | N | - |
| *folC* | *Rv2447c* | PAS | - | 0.66 | 0.78 | 0.11 | Y | N | - |
| *rplC* | *Rv0701* | LZD | - | 0.57 | 0.77 | 0.21 | N | N | - |

743

**Table 3. Homoplastic F<sub>ST</sub> outliers.** Weir and Cockerham's $F_{ST}$ (wcFst) values in the top 1% of values genome wide are reported for each drug. For intergenic SNPs, the closest gene is listed. We identified mutations in genes previously associated with drug resistance (Known = Y) and novel putative resistance or compensatory mutations (Known = N).

| Location | Gene | Type | AMI | CAP | EMB | Et | INH | K | MOX | OFL | PRO | PZA | RIF | STR | Known | Lineage |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1821 | dnaN | intergenic | - | - | - | 0.43 | - | 0.57 | - | 0.10 | - | - | - | - | N | all |
| 7570 | gyrA | missense | - | 0.33 | 0.11 | 0.46 | - | 0.66 | - | 0.29 | - | 0.18 | - | - | Y | all |
| 7572 | gyrA | missense | - | - | - | - | - | - | 0.06 | - | - | - | - | - | Y | all |
| 7581 | gyrA | missense | - | - | - | - | - | - | 0.07 | - | - | - | - | - | Y | all |
| 7582 | gyrA | missense | - | - | - | - | - | - | 0.35 | 0.22 | - | - | - | - | Y | all |
| 75233 | icd2 | intergenic | - | - | - | - | - | - | 0.05 | - | - | - | - | - | N | all |
| 94388 | hycQ | synonymous | - | - | 0.07 | - | 0.12 | - | - | - | - | - | 0.12 | 0.13 | N | all |
| 230170 | Rv0194 | missense | - | - | - | - | 0.12 | - | 0.05 | - | - | - | 0.12 | 0.13 | N | all |
| 332916 | vapC25 | missense | - | - | 0.10 | - | - | - | - | 0.09 | - | 0.20 | - | - | N | all |
| 761155 | rpoB | missense | - | - | 0.31 | - | 0.58 | - | - | - | - | 0.10 | 0.72 | 0.41 | Y | all |
| 761161 | rpoB | missense | - | 0.33 | 0.09 | 0.51 | - | 0.71 | - | 0.13 | - | 0.16 | - | - | Y | all |
| 764817 | rpoC | missense | 0.19 | - | - | - | - | - | - | - | - | - | - | - | N | all |
| 781687 | rpsL | missense | - | - | 0.10 | - | 0.32 | - | - | - | - | 0.15 | - | 0.37 | Y | all |
| 922004 | Rv0830 | missense | - | 0.30 | 0.12 | 0.43 | - | - | - | 0.10 | - | 0.21 | - | - | N | all |
| 1076880 | Rv0965c | synonymous | - | - | - | - | 0.12 | - | - | - | - | - | 0.12 | 0.13 | N | all |
| 1673425 | fabG1 | intergenic | - | - | - | - | - | - | - | - | 0.11 | - | - | - | Y | all |
| 1673432 | fabG1 | intergenic | - | - | - | 0.52 | - | 0.65 | - | - | - | - | - | - | Y | all |
| 1722228 | pks5 | missense | - | - | 0.08 | - | 0.28 | - | - | 0.07 | - | 0.17 | - | 0.26 | N | all |
| 2122395 | lldD2 | synonymous | - | - | - | - | - | - | 0.06 | - | - | - | - | - | N | all |
| 2155168 | katG | missense | - | - | 0.36 | - | 0.89 | - | - | 0.13 | - | 0.32 | 0.60 | 0.66 | Y | all |
| 2174216 | Rv1922 | synonymous | - | - | - | - | - | - | - | - | - | 0.08 | - | - | N | all |
| 2207525 | Rv1958c | intergenic | - | - | - | - | - | - | - | - | - | 0.09 | - | - | N | all |
| 2422824 | Rv2161c | missense | - | 0.30 | - | 0.43 | - | 0.57 | - | 0.10 | - | - | - | - | N | all |
| 2660319 | mbtF | missense | - | - | 0.06 | - | - | - | - | - | - | - | - | - | N | all |
| 2715369 | Rv3413c | intergenic | 0.17 | - | 0.09 | - | 0.28 | - | - | - | - | - | 0.30 | 0.13 | N | all |
| 2866647 | lppA | synonymous | - | - | - | - | 0.12 | - | 0.07 | - | - | - | - | - | N | all |

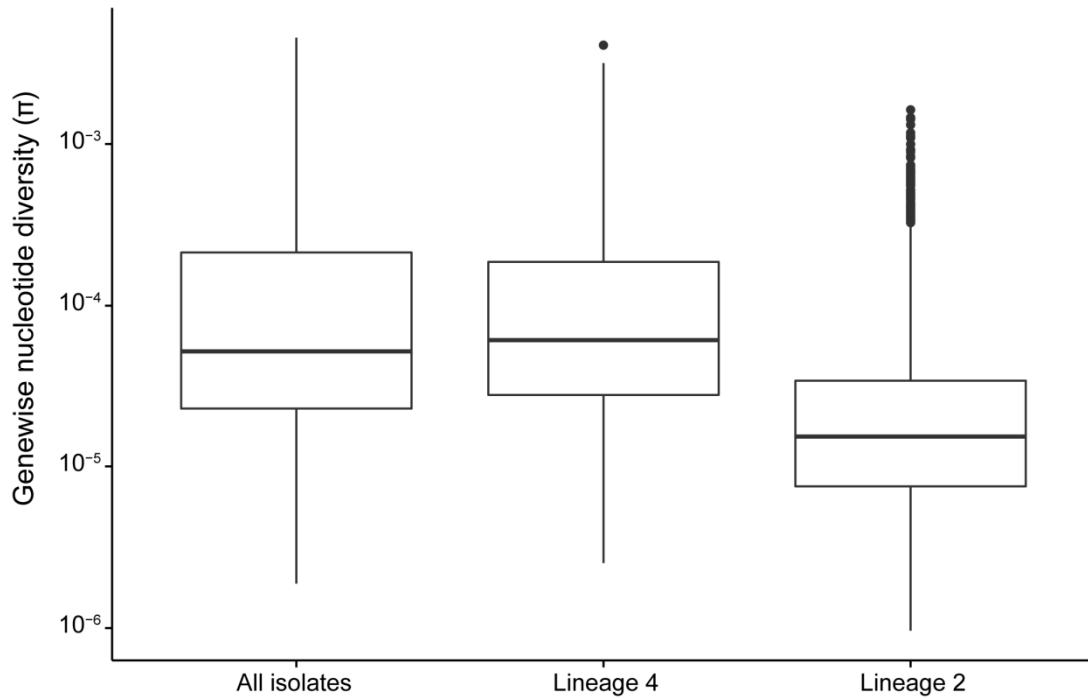| Position | Gene | Type | | | | | | | | | | | | | | |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 2867298 | *lppB* | synonymous | - | - | - | - | 0.13 | - | - | - | - | - | - | - | N | all |
| 2867347 | *lppB* | synonymous | - | - | - | - | 0.13 | - | 0.06 | - | - | - | 0.12 | 0.14 | N | all |
| 2867756 | *lppB* | synonymous | - | - | - | - | 0.14 | - | - | - | - | - | - | - | N | all |
| 3500149 | *Rv3134c* | synonymous | - | - | - | - | - | - | - | - | - | - | 0.11 | - | N | all |
| 3550789 | *Rv3183* | synonymous | - | - | - | - | - | - | - | - | - | - | 0.12 | 0.13 | N | all |
| 3680932 | *lhr* | synonymous | - | - | - | - | 0.12 | - | - | - | - | - | 0.12 | 0.13 | N | all |
| 4001622 | *fadA6* | intergenic | - | - | - | - | - | - | - | - | - | - | 0.11 | - | N | all |
| 4247429 | *embB* | missense | - | - | 0.25 | 0.45 | 0.23 | - | 0.05 | 0.11 | - | 0.31 | 0.21 | 0.20 | Y | all |
| 4247574 | *embB* | synonymous | 0.19 | - | 0.07 | - | 0.27 | - | - | - | - | - | 0.30 | - | Y | all |
| 4327480 | *ethA* | intergenic | 0.20 | - | 0.07 | - | 0.27 | - | - | - | - | - | 0.30 | - | Y | all |
| 764948 | *rpoC* | missense | - | - | - | - | - | - | - | - | 0.06 | - | - | - | Y | L2 |
| 4248003 | *embB* | missense | - | 0.16 | - | - | - | - | - | - | - | - | - | - | Y | L2 |
| 698 | *dnaA* | missense | - | - | - | - | - | - | - | - | - | - | - | 0.10 | N | L4 |
| 60185 | *Rv0057* | missense | - | - | - | - | - | - | - | - | - | - | - | 0.06 | N | L4 |
| 761110 | *rpoB* | missense | 0.66 | - | - | - | - | - | - | - | - | - | - | - | Y | L4 |
| 764822 | *rpoC* | missense | - | - | - | - | - | - | - | - | - | - | - | 0.06 | Y | L4 |
| 781822 | *rpsL* | missense | - | - | - | - | 0.12 | - | - | - | - | - | 0.13 | 0.14 | Y | L4 |
| 2123145 | *lldD2* | missense | - | - | - | - | - | - | - | - | - | - | - | 0.06 | N | L4 |
| 2372550 | *dop* | missense | 0.64 | - | - | - | - | - | - | - | - | - | - | - | N | L4 |
| 2715344 | *Rv2413c* | intergenic | - | - | - | - | - | - | - | - | - | - | - | 0.06 | N | L4 |
| 2986827 | *Rv2670c* | missense | - | - | - | - | 0.16 | - | - | - | - | - | 0.17 | 0.15 | N | L4 |
| 4247431 | *embB* | missense | - | - | - | - | 0.11 | - | - | - | - | - | 0.11 | 0.07 | Y | L4 |
| 4248003 | *embB* | missense | - | - | - | - | - | - | - | - | - | - | - | 0.06 | Y | L4 |

747

748

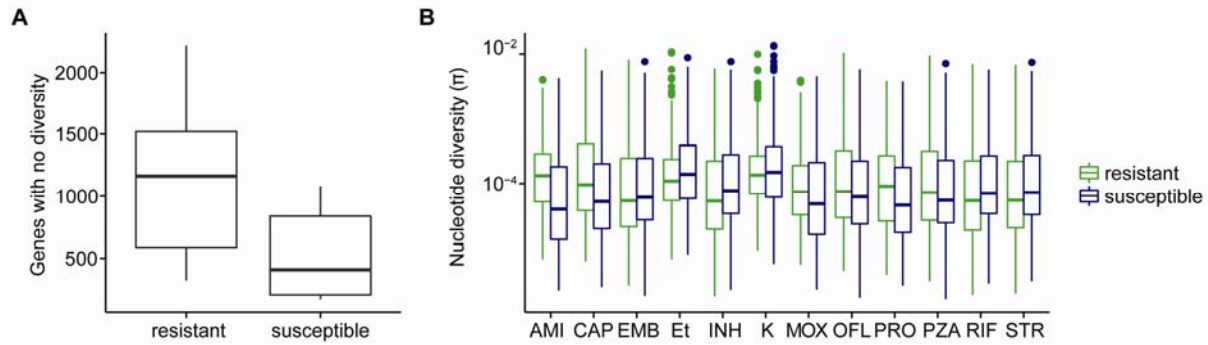749

**Figure 1. Phylogeny of *Mycobacterium tuberculosis* sample.** The phylogeny was inferred using FastTree (60). Lineages are colored as follows: lineage 1 (L1) - pink, lineage 2 (L2) - blue, lineage 3 (L3) - purple, lineage 4 (L4) - red. Lineage 4 is associated with deeper branching sub-lineages in comparison with lineage 2.
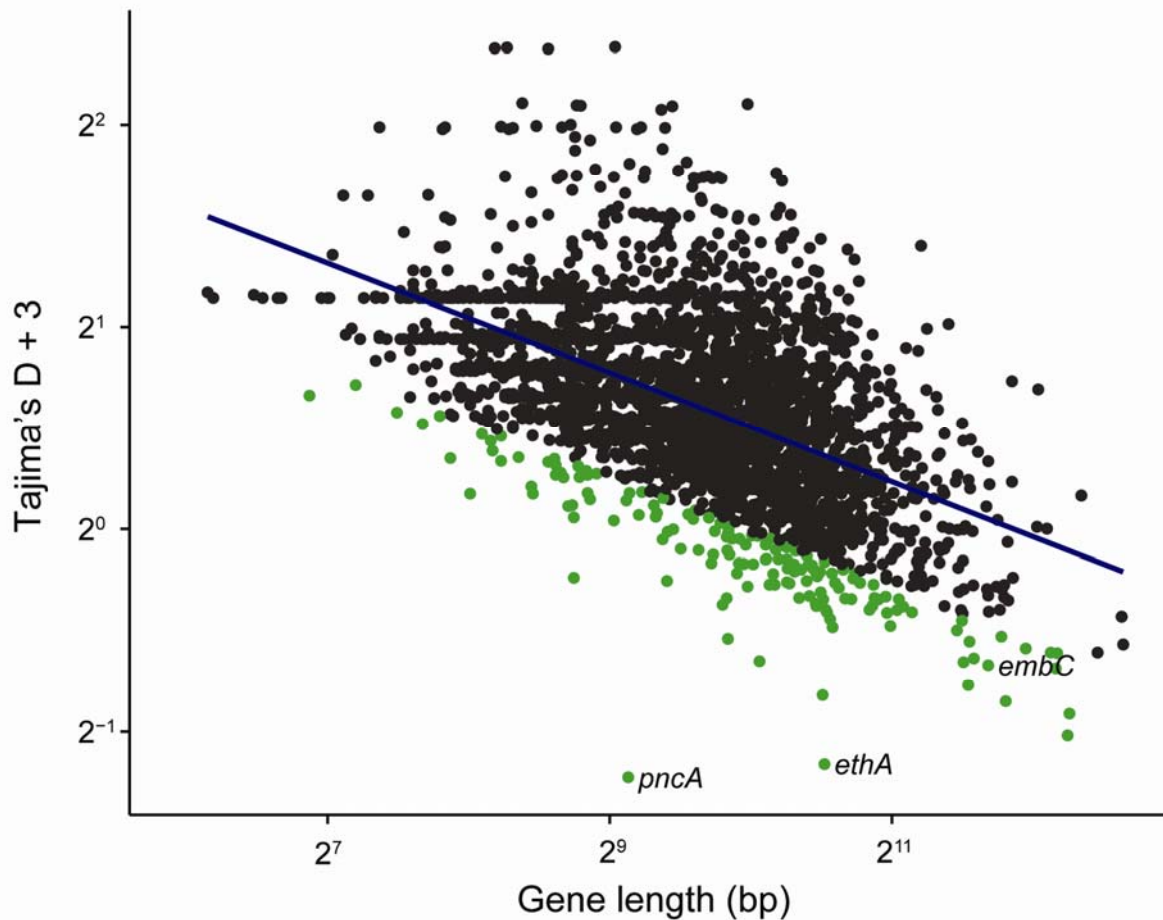
754

**Figure 2. Distributions of gene-wise nucleotide diversity for all isolates, as well as lineages 4 and 2 considered separately.** Repetitive regions of the alignment were masked. Sites were included in estimation of π if 95% of isolates in the alignment had a valid nucleotide at the position. We used Egglib to calculate statistics (64). Nucleotide diversity is lower in lineage 2 compared to lineage 4 (Welch Two Sample t-test, $p < 2.2 \times 10^{-16}$)
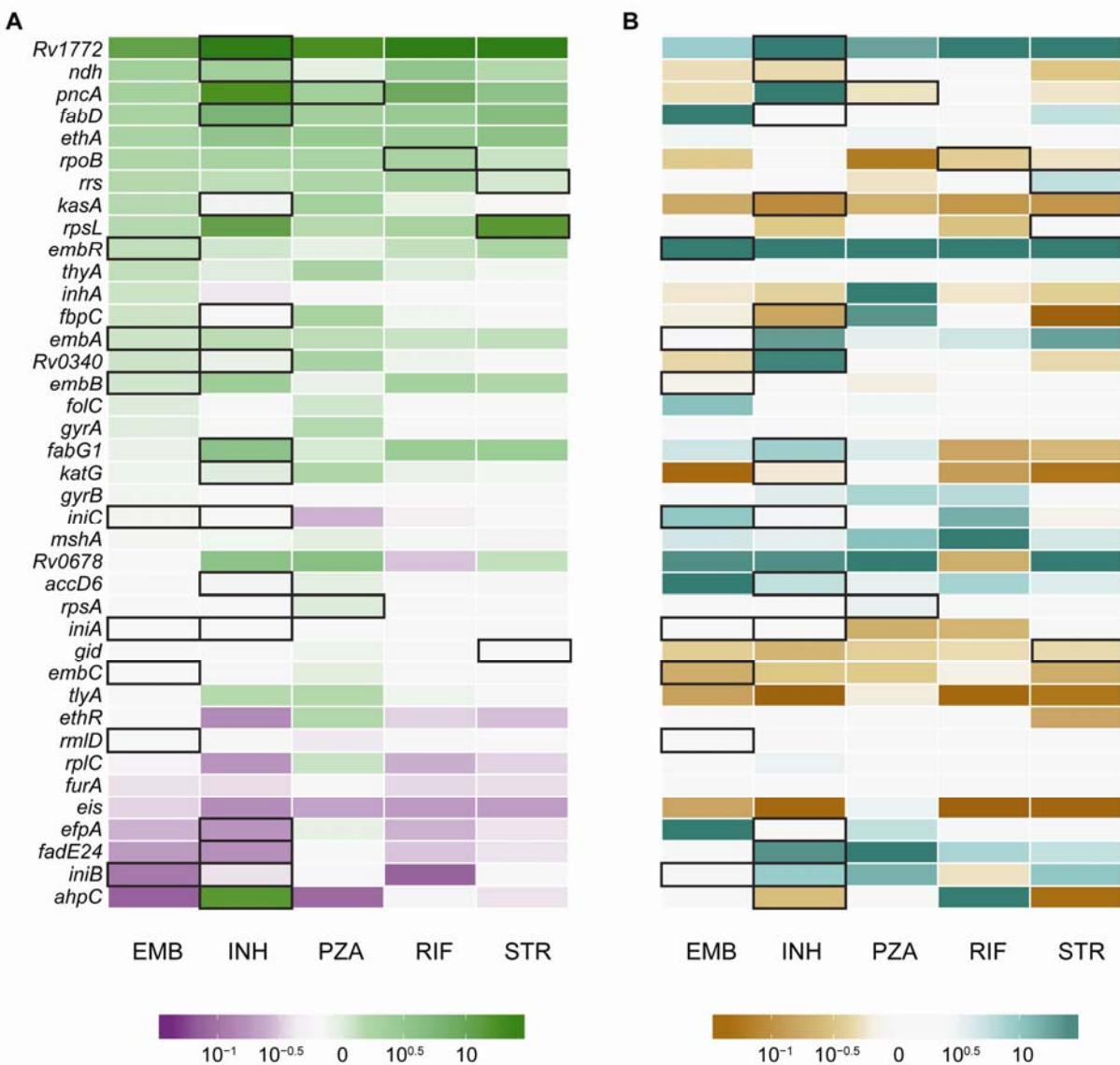
761

762
763 **Figure 3. Diversity of resistant and susceptible isolates.** A) Counts of genes with no
764 nucleotide diversity in resistant and susceptible subpopulations. B) Genewise nucleotide
765 diversity (excluding invariant genes) in susceptible and resistant isolates. Among genes
766 in which it is measurable, nucleotide diversity is similar between resistant and
767 susceptible isolates even when drug resistance associated genes and targets of
768 independent mutation identified by Farhat et al. 2013 are removed ($p = 0.13$).
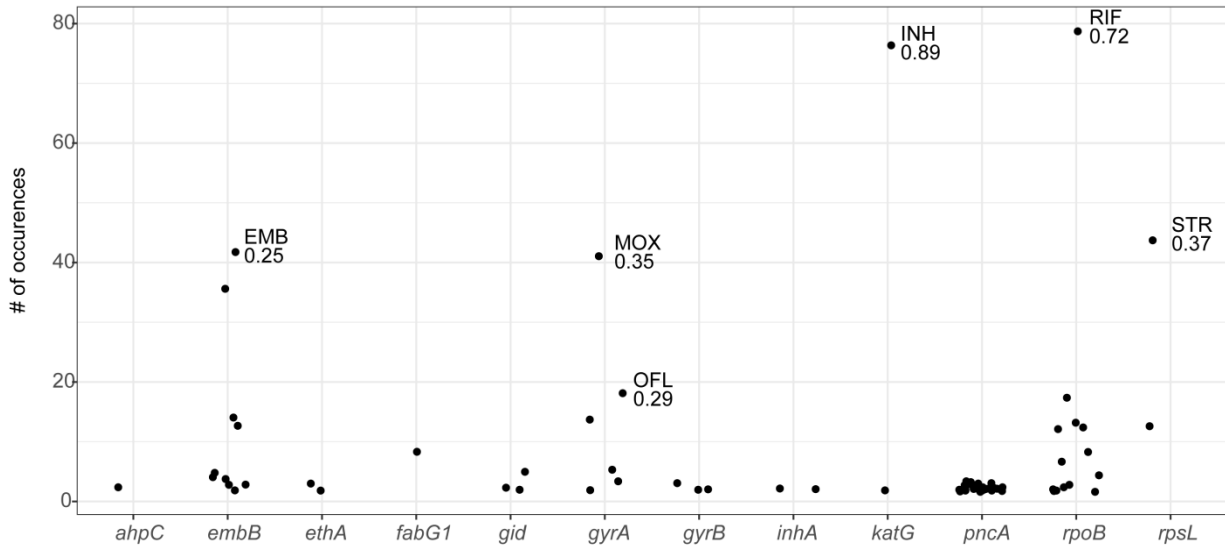
769

**Figure 4. Gene-wise Tajima's D and gene length.** Repetitive regions of the alignment were masked. Gene lengths have been log transformed (base 2). We added a constant value (3) to all Tajima's D values to make them positive and log transformed (base 2), as with the gene lengths. The linear regression line is plotted in blue. Genes with regression values in the lower 5% are highlighted in green. Drug resistance associated genes in this group are labelled. While negative Tajima's D is normally associated with purifying selection or a recent selective sweep, we find that drug resistance genes with negative Tajima's D also have high nucleotide diversity. We hypothesize that patterns of diversity at these genes have been affected by relaxation of purifying selection and positive selection in association with for drug resistance.

**Figure 5. Ratios of nucleotide diversity in resistance associated genes.** Genes with zero diversity were transformed to $1 \times 10^{-16}$ before calculating ratios. Genes with ratios more extreme than $10^{-1.5}$ or $10^{1.5}$ are all filled with the deepest shade. Genes associated with resistance to each drug are outlined in black. A) Ratio of nucleotide diversity in resistant and susceptible isolates. Green genes are more diverse in resistant isolates, which could be due to diversifying selection and/or relaxation of purifying selection. Purple genes are more diverse in susceptible isolates, likely due to increased purifying selection. White genes have similar diversity in resistant and susceptible isolates. B) Comparison of ratios in lineage 2 and lineage 4. Teal genes are more diverse in lineage 2 resistant isolates, suggesting diversifying selection/relaxation of purifying selection specific to this lineage. Brown genes are more diverse in lineage 4 resistant isolates. White genes have similar diversity in lineages 2 and 4.

**Figure 6. Homoplastic SNPs in drug resistance associated genes.** SNPs with $F_{ST}$ in the top 1% of genome-wide values are labeled with the population (associated drug resistance) and the $F_{ST}$ value. *pncA* is remarkable for harboring diverse homoplastic mutations, each of which occurs relatively infrequently ("sloppy target"). *embB*, *gyrA*, *katG*, *rpoB* and *rpsL* harbor dominant mutations that occur frequently on the phylogeny and are strongly associated with resistant populations ("tight targets").

802    **Supplementary Table 1.** Accession numbers and lineage designation for sequence

803    data passing quality control filters.